

# CONTEXT: Convolutional Neural Network and LSTM Code for Text Categorization

May 27, 2016

## Revision history

- May, 2016: Revised when v3.00 (corresponding to [3]) was released. In particular, Sections 4 and 5 were added for LSTM, and the description of memory usage in Section 7 was revised.
- September, 2015: Revised when v2.00 (corresponding to [2]) was released. In particular, Section 3 was added for semi-supervised learning, and description of a *word-mapping file* was added to Section 1.
- August 12, 2015: Section 7 was revised to provide hints and tips on consumption of GPU memory and CPU memory during training and when to make data batches.

## Introduction

This document describes how to use CONTEXT, neural network code for text categorization available at [http://riejohnson.com/cnn\\_download.html](http://riejohnson.com/cnn_download.html).

CONTEXT provides an implementation of *convolutional neural networks (CNNs)* for text categorization described in [1, 2]. These papers describe application of CNN directly to one-hot vectors representing words (*one-hot CNN*), which leads to internally learning an embedding of text regions (or context as in the name) directly from one-hot vectors. [1] considers the end-to-end supervised setting (no additional unlabeled data or no additional algorithm for pre-training), and [2] considers the semi-supervised setting that benefits from unlabeled data as additional resource.

Starting from v3.00, CONTEXT also provides an implementation of *Long Short-Term Memory (LSTM)* networks for text categorization described in [3].

CONTEXT efficiently handles high-dimensional and sparse documents of variable sizes. Although the original purpose of providing the code was to enable reproducing the results in [1, 2, 3], it can also be used for word-vector CNN (which takes word vectors as input), fully-connected neural networks, and more.

In a basic supervised learning setting, training and testing using CONTEXT is done in the following two steps.

1. **Data preparation:** Generate data files used for training and testing as input. In particular, to speed up training, we generate region vectors (vectors that represent small regions of documents) used in the bottom convolution (or LSTM) layer for training and testing, and write to files in advance, instead of making them on the fly during training/testing.
2. **Training and testing:** Train a neural network and optionally test it on development/test data as training proceeds, using the files generated above as input. Also, the trained network can be saved to a file and applied to other data.

Data preparation is described in Section 1, the interface for supervised one-hot CNN is in Section 2, additional functions and parameters for semi-supervised one-hot CNN are introduced in Section 3, supervised and semi-supervised LSTM is described in Sections 4 and 5, respectively, and Section 7 provides hints and tips.

**Prerequisite** Readers are assumed to be familiar with the concepts/terminology used in [1]. [2] is prerequisite for the sections on semi-supervised models (either CNN or LSTM), and [3] is prerequisite for the sections on LSTM.

**NOTE1** To use `CONTEXT`, your system must be equipped with a CUDA-capable GPU such as Tesla K20, and CUDA needs to be installed. `README` should be referred to for more details of hardware/software requirements and how to build executables.

**NOTE2** The following sections refer to sample scripts located in the `sample` directory. The sample scripts are only intended for describing the interface using small data, and their settings are not optimized for performance. The scripts in the `exp-in-papers` directory are for reproducing the experiments in [1, 2, 3]. These scripts should be referred to for realistic examples on the real-world data.

**NOTE3** Binary files produced by `CONTEXT` are endian sensitive and cannot be shared among the systems with different endianness.

## Contents

|   |           |
|---|-----------|
| <b>1 Data preparation: <code>prepText</code></b>  | <b>4</b>  |
| 1.1 <code>prepText gen_vocab</code> : vocabulary file generation                            | 4         |
| 1.2 <code>conText gen_regions</code> : region file and target file generation               | 5         |
| 1.3 Other functions of <code>prepText</code>  | 6         |
| <b>2 One-hot CNN training and testing: <code>conText</code></b>                             | <b>7</b>  |
| 2.1 <code>conText train</code> : single-connection CNN                                      | 8         |
| 2.1.1 Network parameters  | 8         |
| 2.1.2 Layer parameters  | 9         |
| 2.1.3 Output  | 11        |
| 2.2 <code>conText train</code> : multi-connection CNN                                       | 11        |
| 2.3 <code>conText train</code> : data batches for large training data                       | 12        |
| 2.4 <code>conText train</code> : warm-start from a saved CNN                                | 13        |
| 2.5 <code>conText train</code> : AdaDelta   | 13        |
| 2.6 <code>conText train</code> : rmsprop  | 14        |
| 2.7 <code>conText train</code> : layer-by-layer training of single-connection networks      | 14        |
| 2.8 <code>conText predict</code> (applying a CNN to new data)                               | 14        |
| 2.9 <code>conText write_features</code> (write internal features to a file)                 | 15        |
| 2.10 GPU-related settings   | 16        |
| 2.10.1 Device memory handler  | 16        |
| 2.10.2 To control computation on GPU  | 16        |
| <b>3 Semi-supervised one-hot CNN (using <code>prepText</code> and <code>conText</code>)</b> | <b>17</b> |
| 3.1 Training with unlabeled data (tv-embedding learning)                                    | 17        |
| 3.1.1 Training data generation from unlabeled data  | 18        |
| 3.1.2 Training of tv-embedding  | 19        |
| 3.2 Training with labeled data and tv-embedding   | 20        |
| <b>4 LSTM training and testing: <code>reText</code></b>                                     | <b>20</b> |
| 4.1 <code>reText train</code> : training and testing LSTM                                   | 21        |
| 4.1.1 Network parameters  | 21        |
| 4.1.2 Layer parameters  | 21        |
| 4.1.3 Output  | 22        |
| 4.2 <code>reText predict</code> (applying a saved model to new data)                        | 22        |

|          |   |           |
|----------|---|-----------|
| <b>5</b> | <b>Semi-supervised LSTM (using <code>prepText</code> and <code>reText</code>)</b>   | <b>23</b> |
| 5.1      | LSTM training with unlabeled data (LSTM tv-embedding learning)                      | 23        |
| 5.1.1    | Training data generation from unlabeled data  | 23        |
| 5.1.2    | Training of LSTM tv-embedding   | 24        |
| 5.2      | Training with labeled data and LSTM tv-embedding                                    | 24        |
| <b>6</b> | <b>Examples of other usage</b>  | <b>25</b> |
| 6.1      | One-hot CNN with a convolution layer above a convolution layer                      | 25        |
| 6.2      | One-hot CNN with a fully-connected layer above a convolution layer                  | 25        |
| 6.3      | Fully-connected neural networks   | 25        |
| 6.4      | Linear models   | 25        |
| 6.5      | Word-vector CNN   | 25        |
| <b>7</b> | <b>Hints and Tips</b>   | <b>26</b> |
| 7.1      | Miscellaneous tips  | 26        |
| 7.2      | GPU memory consumption during training ( <code>conText train, reText train</code> ) | 27        |
| 7.3      | CPU memory consumption during training ( <code>conText train, reText train</code> ) | 27        |
| 7.4      | When to make training data batches ( <code>conText train, reText train</code> )     | 28        |

# 1 Data preparation: prepText

*To get started* Try `sample/sample.sh`.

An executable `prepText` prepares input data for training and testing. It runs on CPU and does not use GPU.

Usage: `prepText action param1 param2 ...`

The first argument *action* must be one of the following:

- `gen_vocab`: Generate a vocabulary file used as input to `gen_regions` and other functions of `prepText`.
- `gen_regions`: Generate a region vector file, a target file, and a word-mapping file used as input to CNN training/testing and LSTM training/testing.
- `show_regions`: Show the content of a region vector file.
- `gen_nbw`: Generate a NB weight file, used as input to `gen_nbwfeat`.
- `gen_nbwfeat`: Generate a NB-weighted feature file, used in the seq2-bow $n$ -CNN experiments in [1].
- `gen_b_feat`: Generate a bag-of- $n$ -gram file, used in the baseline experiments with SVM and fully-connected neural networks in [1].

The rest of the arguments depend on the *action* and they are described below.

## 1.1 prepText gen\_vocab: vocabulary file generation

As the first step of data preparation, we generate a vocabulary file from training data. That is, we extract words (or  $n$ -grams) from text data.

**Tokenized text file (input)** Text files given to `gen_vocab` as input should contain documents, one per line, and each document should be already tokenized so that tokens are delimited by space. See `sample/data/s-test.tok` for an example.

**Vocabulary file (output)** The vocabulary file generated by `gen_vocab` contains words (or  $n$ -grams), one per line. Each word is optionally followed by its frequency; in that case, a tab character is used as a delimiter between the word and the frequency.

**To reduce vocabulary size** Typically, it is not necessary to use all the words in the training data to obtain good text categorization performance, while retaining many words is resource consuming. `gen_vocab` provides several ways to reduce the size of vocabulary such as only retaining the most frequent words up to `max_vocab_size`, only retaining the words occurring no fewer than `min_word_count`, removing stopwords listed in `stopword_fn`, and removing the words that include numerical characters (`RemoveNumbers`).

In the tables below, required parameters are marked with \*; all other parameters are optional.

---

| Parameters for <code>prepText gen_vocab</code> (vocabulary file generation). |  |
|--|--|
| * <code>input_fn=</code>   | Path to the input token file.  |
| * <code>vocab_fn=</code>   | Path to the file that the vocabulary is written to.  |
| <code>stopword_fn=</code>  | Path to a stopwords file. The words in this file will be excluded.   |
| <code>min_word_count=</code>   | Minimum word counts to be included in the vocabulary file. (Default:No limit)  |
| <code>max_vocab_size=</code>   | Maximum number of words to be included in the vocabulary file. The most frequent ones will be included. (Default:No limit) |
| <code>LowerCase</code>   | Convert upper-case to lower-case characters.   |

|               |  |
|---------------|--|
| UTF8          | Convert UTF8 en dash, em dash, single/double quotes to ascii characters.         |
| RemoveNumbers | Exclude words that contain numbers.  |
| WriteCount    | Write word counts as well as the words to the vocabulary file.                   |
| n=            | <i>n</i> for <i>n</i> -grams. E.g., if <i>n</i> =3, only tri-grams are included. |

---

**Example 1** (`sample/sample.sh`). *Extract words (delimited by space) from `data/s-train.tok` and convert them to lower cases while only retaining the most frequent 10000 words:*

```
prepText gen_vocab input_fn=data/s-train.tok vocab_fn=data/s.vocab \
LowerCase max_vocab_size=10000
```

## 1.2 conText gen\_regions: region file and target file generation

Next, we generate a region file and target file, which serve as input files for CNN (or LSTM) training and testing.

**Tokenized text file (input)** The input text file format is the same as `gen_vocab` above. The file should contain one document per line, and each document should be already tokenized so that tokens are delimited by space.

**Label file (input)** Each line of a label file should contain classification labels for each document, and the order of the documents must be the same as the text file. In case of multi-label classification (i.e., more than one label can be assigned to each document), the labels should be delimited by a vertical line `|`. The labels (any string) must be declared in a label dictionary described below.

**Text/label file naming conventions** The text file and the corresponding label file must have the same pathname stem with different file extensions, e.g., text file `data/s-train.tok` and label file `data/s-train.cat`.

**Label dictionary file (input)** The labels used in the label file above must be declared in a label dictionary file. The label dictionary file should contain one label per line. See `sample/data/s-cat.dic` for example.

**Vocabulary file (input)** The vocabulary file should be generated by `gen_vocab` above. To use a vocabulary file generated by some other means, note that a tab character (0x09) is regarded as a delimiter, i.e., in each line, a tab character and anything that follows are ignored. Also note that the case option must be consistent; e.g., if `gen_region` specifies `LowerCase` (which converts upper-case letters to lower-case), then the contents of the vocabulary file must also be all lower-case.

**Region file (output)** `gen_regions` generates a region file that contains region vectors with the specified region size, stride, and padding, in binary format. To generate the region file, a file extension `.xsmatvar` (which indicates that the file is in a sparse and variable-sized format) is automatically attached to the pathname stem specified by `region_fn_stem=`.

**Target file (output)** `gen_regions` generates a target file that contains classification targets in the format that the executable for CNN training can read. To generate the target file, a file extension either `.y` or `.ysmat` (which indicates the file format) is automatically attached to the pathname stem specified by `region_fn_stem=`.

**Word-mapping file (output)** `gen_regions` generates a *word-mapping file* that contains a mapping between words (or *n*-grams) and dimensions of the region vectors. To generate the word-mapping file, a file extension `.xtext` is automatically attached to the pathname stem specified by `region_fn_stem=`. This file is used by v2.00 (or higher) to ensure the consistency of the mapping at the time of training and application of the models.

**Multi-label/no-label documents** By default, single-label classification is assumed, and therefore, each document is expected to have exactly one label; the process terminates with an error on documents with more than one label or no label. To allow multi-label and no-label documents, specify `MultiLabel`.

---

| <b>Parameters for <code>prepText gen_regions</code> (region file generation).</b> |   |
|---|---|
| * <code>input_fn=</code>  | Input filename without extension.   |
| * <code>text_fn_ext=</code>   | Filename extension of the tokenized text file.  |
| * <code>label_fn_ext=</code>  | Filename extension of the label file. Required for target file generation.  |
| * <code>label_dic_fn=</code>  | Path to the label dictionary file (input). The file should list the labels used in the label files, one per each line. Required for target file generation.   |
| * <code>vocab_fn=</code>  | Path to the vocabulary file generated by <code>gen_vocab</code> (input).  |
| * <code>region_fn_stem=</code>  | Pathname stem of the region vector file, target file, and word-mapping file (output). To make the pathnames of these files, the respective extensions will be attached. The name must not contain a plus sign <code>+</code> .                        |
| <code>y_ext=</code>   | Filename extension of the target file (output). <code>.y</code>   <code>.ysmat</code> . Use <code>.ysmat</code> (binary sparse format) if the number of classes is large. (Default: <code>.y</code> (text format))                                    |
| * <code>patch_size=</code>  | Region size. For LSTM, set this to 1.   |
| <code>patch_stride=</code>  | Region stride. (Default: 1)   |
| <code>padding=</code>   | Padding size. For CNNs, region size minus one is recommended. (Default: 0)  |
| <code>Bow</code>  | Generate region vectors for a bow-convolutional layer. See Section 2.2.2 of [1] for bow-convolution. Shortened from <code>Bow-convolution</code> , which still works.   |
| <code>VariableStride</code>   | Take variable strides; see Section 3.1 of [1].  |
| <code>LowerCase</code>  | Convert upper-case to lower-case characters. On/off of this switch must be consistent with <code>gen_vocab</code> .   |
| <code>UTF8</code>   | Convert UTF8 en dash, em dash, single/double quotes to ascii characters. On/off of this switch must be consistent with <code>gen_vocab</code> .   |
| <code>MultiLabel</code>   | Allow multiple labels per data point for multi-label classification. Data points that do not have any label are also allowed.   |
| <code>AllowZeroRegion</code>  | Do not ignore empty regions.  |
| <code>RegionOnly</code>   | Generate a region file only. Do not generate a target file.   |
| <code>batch_id=</code>  | Batch ID, e.g., <code>1of5</code> (the first batch out of 5), <code>2of5</code> (the second batch out of 5). Specify this when making multiple sets of files ( <i>batches</i> ) for one large dataset. See Section 2.3 for how to use this parameter. |

---

**Example 2** (`sample/sample.sh`). *Generate a region vector file `data/s-train-p3.xsmatvar`, a target file `data/s-train-p3.y`, and a word-mapping file `data/s-train-p3.xtext` from a tokenized text file `data/s-train.tok` and a label file `data/s-train.cat`, with region size 3, stride 1, and padding 2.*

```
prepText gen_regions input_fn=data/s-train text_fn_ext=.tok label_fn_ext=.cat \
                    vocab_fn=s.vocab LowerCase \
                    label_dic_fn=data/s-cat.dic \
                    patch_size=3 patch_stride=1 padding=2 \
                    region_fn_stem=data/s-train-p3
```

### 1.3 Other functions of `prepText`

The parameters for the other functions are similar to those described above. To display help, enter

```

prepText  show_regions  or
prepText  gen_nbw       or
prepText  gen_nbwfeat   or
prepText  gen_b_feat

```

For example usage, see the scripts `exp-in-papers/train_imdb_seq2_bown.sh` for `gen_nbw` and `gen_nbwfeat`, and `sample/sample-1nn.sh` for `gen_b_feat`.

## 2 One-hot CNN training and testing: `conText`

**To get started** Try `sample/sample.sh` or `exp-in-papers/train_imdb_seq.sh`.

`conText` is an executable for training and testing (mainly) CNNs using GPU.

Usage: `conText gpu#[:memory_size] action param1 param2 ...`

**gpu#** The first argument specifies the GPU ID number. Specify `-1` if the default GPU should be used,

**memory\_size** Description of this argument is delayed to Section 2.10.1.

**action** The second argument *action* must be one of the following:

- `train`: Train a CNN. Optionally, test the CNN on development/test data as training proceeds. This function can also be used to train a fully-connected neural network or a linear model.
- `predict`: Apply the trained model to new data.
- `write_features`: Write internal features to a file.

**param1, param2, ...** Parameters can be specified via either command line arguments or a parameter file or a combination of both. If an argument starts with an `@`, the string following `@` is regarded as a pathname to a parameter file. See `sample/param/sample.param` for an example of the parameter file. Use of a parameter file is sometimes convenient as parameters can be organized and commented.

**Example 3** (`sample/sample-paramfile.sh`). Use parameters in `param/sample.param` and set `random.seed=1`.

```
conText -1 train random.seed=1 @param/sample.param
```

To describe the interface, the following term/concepts will be used.

- **Hidden layer**: In `conText`, each *hidden layer* performs a sequence of operations: weighting and activation (computation of  $\sigma(\mathbf{W}\mathbf{x}+\mathbf{b})$ ), pooling (optional), and response normalization (optional), in that order. (Note that an alternative is to treat each function as an individual layer, e.g., a pooling layer and a response normalization layer, and so on.) We assign unique id numbers to hidden layers starting from zero, e.g., layer-0, layer-1, ...
- **Top layer**: A top layer is a special layer that only performs weighting ( $\mathbf{W}\mathbf{x} + \mathbf{b}$ ) to generate output.
- **Single-connection architecture**: With a *single-connection* architecture, data flows from layer-0 to layer-1, from layer-1 to layer2, and so on and from the last hidden layer to the top layer, as illustrated in Figure 1.
- **Multi-connection architecture**: A *multi-connection* architecture allows each layer to be connected to more than one layer; e.g., see Figure 2.

First, Section 2.1 describes the interface for training a CNN with single-connection architecture. CNN with a multi-connection architecture only requires a few additional parameters, described later.

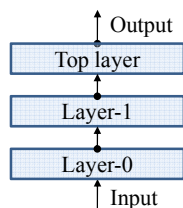


Figure 1: Single-connection example.

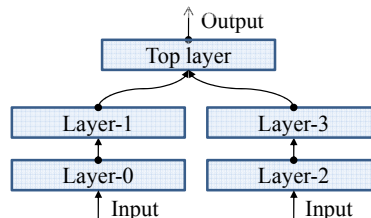


Figure 2: Multi-connection example.

## 2.1 conText train: single-connection CNN

**To get started** Try `sample/sample.sh` or `exp-in-papers/train_imdb_seq.sh`.

There are two types of parameters, *layer parameters* which specifies parameters for each layer, and *network parameters* which are not specific to layers.

### 2.1.1 Network parameters

The parameters below are for controlling the overall (as opposed to layer-specific) aspect of training.

| <b>To control input</b>                |   |
|--|---|
| * trnname=                             | Training data name. Region/target filename (Section 1.2) without extension.   |
| * tstname=                             | Test/development data name. Region/target filename (Section 1.2) without extension. It must not contain a plus sign. Required for testing unless NoTest is specified.   |
| * datatype=sparse                      | Dataset type. Sparse input. Other options are described later.  |
| data_dir=                              | Directory where data files are.   |
| * x_ext=.xsmatvar                      | Region filename extension. Other types are described later.   |
| * y_ext=                               | Target filename extension. <code>.y</code> (text format)   <code>.ysmat.</code> (sparse format)   |
| NoTest                                 | Do no testing. Training only.   |
| <b>To control output</b>               |   |
| save_fn=                               | Path to the file the model is saved to for warm-start later.  |
| save_interval=                         | <i>s</i> . The model is saved to a file every <i>s</i> epochs. (Default: no save)   |
| test_interval=                         | <i>t</i> . Evaluate performance on test/development data every <i>t</i> epochs. (Default: 1)  |
| evaluation_fn=                         | Write performance results to this file in the csv format.   |
| ExactTrainingLoss                      | Compute and display training loss. If this switch is off, estimated loss is shown instead.  |
| LessVerbose                            | Do not display too much information.  |
| inc=                                   | <i>d</i> . Show progress after going through <i>d</i> data points. Useful when data is large or when training is time-consuming.  |
| <b>To define network architecture.</b> |   |
| * layers=                              | Number of hidden layers.  |
| <b>To control training.</b>            |   |
| * num.iterations=                      | Number of epochs.   |
| * loss=                                | Loss function. Square   Log   BinLogi2<br>Square: $(p - y)^2/2$ where <i>p</i> is prediction and <i>y</i> is target.<br>Log (softmax and log loss)<br>BinLogi2 (binary logistic loss for $y \in \{0, 1\}$ ): $\log(1 + \exp(-(2y - 1)p))$ |
| max_loss=                              | When estimated training loss goes over this value, training is terminated. Use this to detect 'explosion' typically caused by setting the learning rate to a too large value. Only available with v3.00 or higher.                        |



|                                   |   |
|-----------------------------------|---|
| <code>mini_batch_size=</code>     | Mini-batch size for training. (Default:100)   |
| <code>step_size_scheduler=</code> | Step-size scheduler type. Few (reduce it a few times) is recommended. (Default: no scheduling)  |
| <code>step_size_decay=</code>     | Use this with <code>step_size_scheduler=Few</code> . Step-size is reduced by multiplying this value when it is reduced.   |
| <code>step_size_decay_at=</code>  | Use this with <code>step_size_scheduler=Few</code> . Step-size is reduced after this many epochs. To reduce it more than once, use an underbar ( <code>_</code> ) as a delimiter, e.g., <code>80_90</code> reduces it after the 80-th and 90-th epochs. |

---

#### Other parameters.

|                                    |  |
|------------------------------------|--|
| <code>random_seed=</code>          | Seed for random number generation.   |
| <code>Regression</code>            | Specify this if the task is regression. (Default: classification)  |
| <code>MultiLabel</code>            | Specify this if each data point can be assigned multiple labels. (Default: single-label classification). |
| <code>test_mini_batch_size=</code> | Mini-batch size for parallel processing for testing. (Default: 100).                                     |

---

**Example 4** (`sample/sample.sh`). *Perform 20 epochs of SGD with mini-batch size 100 while reducing the step-size by multiplying 0.1 after 15 epochs.*

```
... num_iterations=20 mini_batch_size=100 \
    step_size_scheduler=Few step_size_decay_at=15 step_size_decay=0.1
```

### 2.1.2 Layer parameters

The parameters below can be specified for a specific layer by attaching the layer id  $\ell$  ( $\ell = 0, 1, \dots, \text{top}_-$ ) in front of the keywords, e.g., `0step_size=0.1` for layer-0 or `top_reg_L2=1e-4` for the top layer. Without a layer id, the parameter is regarded as a default parameter and applied to all the layers; i.e., the parameter with a layer id supersedes the one without the id. For example,

```
layers=3 reg_L2=0.0001 0reg_L2=0 top_reg_L2=0.01
```

is equivalent to

```
layers=3 0reg_L2=0 1reg_L2=0.0001 2reg_L2=0.0001 top_reg_L2=0.01
```

---

#### Weighting and activation.

|   |  |
|---|--|
| * <code>[<math>\ell</math>]nodes=</code>    | Number of weight vectors (or neurons).                               |
| <code>[<math>\ell</math>]activ_type=</code> | Activation type. None   Log   Rect   Softplus   Tanh. (Default:None) |
|   | None: $\sigma(x) = x$  |
|   | Log (sigmoid): $\sigma(x) = 1/(1 + \exp(-x))$                        |
|   | Rect (rectifier): $\sigma(x) = \max(x, 0)$                           |
|   | Softplus: $\sigma(x) = \log(1 + \exp(x))$                            |
|   | Tanh: $\sigma(x) = (\exp(2x) - 1)/(\exp(2x) + 1)$                    |

---

#### Pooling.

|   |  |
|---|--|
| <code>[<math>\ell</math>]pooling_type=</code>   | Pooling type. Max   Avg   L2   None. (Default:None (no pooling)). To perform pooling, either <code>num_pooling</code> or <code>pooling_size</code> is required, but not both.  |
| <code>[<math>\ell</math>]num_pooling=</code>    | Number of pooling regions. Fixed-sized output is produced; see Section 2.2.3 of [1]. For example, <code>num_pooling=5</code> makes five pooling regions of equal size without overlapping, which leads to output of size 5 (5 ‘pixels’).       |
| <code>[<math>\ell</math>]pooling_size=</code>   | Pooling size, as in CNN for images. Variable-sized output is produced if input is variable-sized. Use this with care since a fully-connected layer such as the top layer requires fixed-sized input. Example: <code>sample/sample-cc.sh</code> |
| <code>[<math>\ell</math>]pooling_stride=</code> | Pooling stride. Required if <code>pooling_size</code> is specified.  |

---

| <b>To control training.</b>         |  |
|-------------------------------------|--|
| <code>[ℓ]dropout=</code>            | Dropout rate [5]. (Default:No dropout). Dropout is applied to the input to the layer. Dropout cannot be applied to sparse input.   |
| <code>[ℓ]reg_L2=</code>             | L2 regularization parameter. Useful value: 1e-4, 1e-3, maybe 1e-5  |
| <code>[ℓ]reg_L2const=</code>        | $c$ : max-norm regularization parameter. Linearly scale to $c$ when the 2-norm of a weight vector exceeds $c$ . <code>reg_L2</code> and <code>reg_L2const</code> are mutually exclusive.   |
| <code>[ℓ]init_weight=</code>        | $x$ : scale of weight initialization. Weights will be initialized by Gaussian distribution with zero mean with standard deviation $x$ . If <code>InitWeightUniform</code> is on, initial weights will be in $[-x, x]$ . (Default:0.01) |
| <code>[ℓ]InitWeightUniform</code>   | Initialize weights with uniform distribution. (Default:Gaussian distribution)  |
| <code>[ℓ]init_intercept=</code>     | Initial values of intercepts. (Default:0)  |
| <code>[ℓ]NoIntercept</code>         | No intercept.  |
| <code>[ℓ]RegularizeIntercept</code> | Regularize intercepts. (Default: intercepts are unregularized)   |
| <b>SGD parameters.</b>              |  |
| * <code>[ℓ]step_size=</code>        | Initial step-size (learning rate) for SGD.   |
| <code>[ℓ]momentum=</code>           | Momentum for SGD. Useful value: 0.9.   |

**Example 5** (`sample/sample.sh`). *In the following, layer-0 is with 500 weight vectors, rectifier activation, and max-pooling with one pooling unit; in all layers, the L2 regularization parameter is 1e-4 and the initial step-size is 0.25; and the top layer performs dropout with rate 0.5.*

```
... 0nodes=500 0pooling_type=Max 0num_pooling=1 0activ_type=Rect \
    reg_L2=1e-4 step_size=0.25 top_dropout=0.5
```

**Response normalization** Optionally, response normalization can be applied, similar to [4]. The pooling result indexed by  $j$  is multiplied with  $(\gamma + \alpha \sum_i \mathbf{v}_i^2)^{-\beta}$  where  $i$  moves within the window of size  $w$  surrounding  $j$ . If `resnorm_type=Cross` (cross-neuron), the sum is taken over neurons in a window defined on a ring (where the first neuron is considered to be next to the last neuron), and in particular, if width  $w$  is no smaller than the number of neurons, the sum is taken over all the neurons.

|                                |  |
|--------------------------------|--|
| <code>[ℓ]resnorm_type=</code>  | None   Cross. (Default:None (do nothing))  |
| <code>[ℓ]resnorm_width=</code> | $w$ : Width $w$ above.                     |
| <code>[ℓ]resnorm_one=</code>   | $\gamma$ in the formula above. (Default:1) |
| <code>[ℓ]resnorm_alpha=</code> | $\alpha$ in the formula above.             |
| <code>[ℓ]resnorm_beta=</code>  | $\beta$ in the formula above.              |

**Region size, stride, and padding with dense input** When `datatype=sparse` or `sparse_multi`, the input files are required to be in the special format that makes use of sparseness (i.e., most vector components are zero), called *sparse input*. With `datatype=sparse` or `sparse_multi`, the first layer takes *sparse input*. Other layers all take *dense input*. With *sparse input*, region size/stride and padding are determined at the time of region file generation using `prepText gen_regions`, not by layer parameters. Region size/stride and padding in the layer with *dense input* are specified via layer parameters as follows.

|                               |   |
|-------------------------------|---|
| <code>[ℓ]patch_size=</code>   | Region size in the convolution layer with <i>dense input</i> .              |
| <code>[ℓ]patch_stride=</code> | Region stride in the convolution layer with <i>dense input</i> .            |
| <code>[ℓ]padding=</code>      | Padding size at the edge in the convolution layer with <i>dense input</i> . |

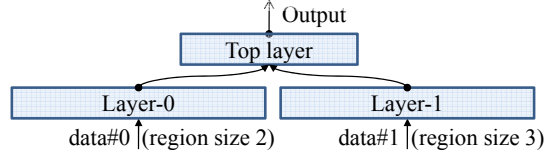


Figure 3: Multi-connection architecture defined in `sample/sample-multiconn.sh`.

If input is fixed-sized, and if `patch_size` is omitted, the layer is considered to be fully-connected. An example of a fully-connected network will be discussed in Section 6

### 2.1.3 Output

**Evaluation results (single-label classification)** Optionally, classification evaluation results can be written in the CSV format to a file specified via `evaluation_fn`. Lines of this file look as follows:

```
ite,5,0.226053, test-loss,0.217621, perf:err,0.34
ite,10,0.119096, test-loss,0.164462, perf:err,0.22
```

The numbers in each line above are the number of epochs, training loss approximation, test loss, and error rate (1 - accuracy). The training loss approximation is the average of data point loss, which is obtained on each data point as a side product of forward propagation during the designated epoch. It is approximate since the network state keeps changing during the epoch. The regularization term is not included. To obtain the real training loss at the expense of slightly longer processing, specify `ExactTrainingLoss`, and it will be computed and shown in the evaluation file, for example, as follows.

```
ite,5,0.226053, trn-loss,0.20709,0.0725723, test-loss,0.217621, perf:err,0.34
```

There are two numbers following `trn-loss`. The first number is the training loss (averaged over data points) excluding the regularization term. The second number is the L2 regularization term.

**Evaluation results (multi-label classification)** The description above was for single-label classification cases. If the task is multi-label classification, indicated by `MultiLabel`, error rate is omitted and only test loss is shown.

**stdout** While evaluation is done in the interval specified by `test_interval`, the loss approximation is written to `stdout` every epoch. If, for example, `inc=100`, the loss approximation is written to `stdout` also every time training goes through 100 more data points, which is convenient with large data. Also, the network and layer setting is written to `stdout` at the beginning of training.

**Model files** When `save_interval` and `save_fn` are specified, models are written to files in binary format. The model files can be used for warm-start of training and making predictions on new data.

## 2.2 conText train: multi-connection CNN

**To get started** Try `sample/sample-multiconn.sh` or `exp-in-papers/train_imdb_seq2.sh`.

Training of a multi-connection CNN requires additional parameters to specify the connection between layers and between input data and layers. We use the parameter setting in the script `sample/sample-multiconn.sh` as a running example. The network architecture of this CNN is shown in Figure 3. The multi-connection parameters in this script are as follows:

```
extension=multi layers=2 conn0=0-top conn1=1-top
datatype=sparse_multi data_ext0=p2 data_ext1=p3
0dataset_no=0 1dataset_no=1
```

and they are explained below.

**To enable multi-connection architecture** `extension=multi` must be specified.

**To define connection** The connections between layers in Figure 3 are defined by:

```
layers=2 conn0=0-top conn1=1-top
```

`layers=2` indicates that there are two hidden layers (as in the single-connection parameter). The numbers on the right to the equal sign are layer id's, and `top` indicates the top layer. `conn0=0-top` indicates that the data flows from layer-0 to the top layer, and `conn1=1-top` indicates that the data flows from layer-1 to the top layer. The numbers following `conn` must be sequential starting from 0.

**Example 6.** *The order of the `conn` parameters does not matter. The connection in Figure 2 can be specified by:*

```
layers=4 conn0=0-1 conn1=2-3 conn2=1-top conn3=3-top or
layers=4 conn0=0-1 conn1=1-top conn2=2-3 conn3=3-top
```

*Although both specifies the same network connection, the results may differ since the order of random initialization of weights may be different.*

**To use multiple datasets** A multi-connection network may have more than one bottom layer. Since regions of the bottom layers are defined at the time of region file generation (`prepText gen_regions`), in order to allow multiple bottom layers to have different region definitions from each other, we need to be able to specify multiple region files as input to training. This is done by first enabling multi-dataset option by:

```
datatype=sparse_multi
```

and then specifying the filename stems of the datasets, in this example, as follows:

```
data_dir=data trnname=s-train- tstname=s-test- data_ext0=p2 data_ext1=p3
```

This indicates that the pathname stems of the two datasets are as follows.

|           | for training    | for testing    |
|-----------|-----------------|----------------|
| dataset-0 | data/s-train-p2 | data/s-test-p2 |
| dataset-1 | data/s-train-p3 | data/s-test-p3 |

Note that actual pathnames to region/target/word-mapping files are obtained by further concatenating the file extensions such as `.xsmatvar` and `.y`. Finally, we specify which layer takes which data as input in the format of `ldataset_no=d`, which indicates layer- $l$  takes dataset- $d$  as input. For example,

```
0dataset_no=0 1dataset_no=1
```

specifies that layer-0 takes dataset-0 and layer-1 takes dataset-1 as input.

### 2.3 conText train: data batches for large training data

When training data is large, it is convenient to divide it to several *batches* (see Section 7.4 for when to use batches). This can be done by doing the following:

1. Divide a tokenized text file and a label file for training data into  $k$  batches (e.g.,  $k = 5$ ).

2. After generating a vocabulary file using `prepText gen_vocab` as before, invoke `prepText gen_regions`  $k$  times to generate a region and a target file for each of  $k$  batches, using `batch_id` to indicate the position in the batches. `prepText gen_regions` uses this parameter to name the region/target files appropriately.
3. Train a CNN with `num_batches=k` to indicate that the training data consists of  $k$  batches. Note that in the case of multiple batches, going through one batch is counted as *one epoch*; therefore, for example, `num_batches=5 num_iterations=100` results in going through the entire data (all five batches) 20 times instead of 100, and with `step_size_decay_at=4`, the step size is reduced after going through about four fifth of the data.

**Example 7** (`sample/sample-batch.sh`). *Example of two batches.*

```

prepText gen_regions  input_fn=s-train1 batch_id=1of2 \
                      region_fn_stem=sb_train ...
prepText gen_regions  input_fn=s-train2 batch_id=2of2 \
                      region_fn_stem=sb_train ...
conText gpu# train  num_batches=2 trnname=sb-train ...

```

## 2.4 conText train: warm-start from a saved CNN

Training can be resumed by warm-starting from a CNN saved to a file. To save a CNN to a file, `save_fn` and `save_interval` should be specified as described 2.1. With `fn_for_warmstart` specified, training warm-starts from the model saved in the designated file.

---

### To warm-start from a saved model

---

|   |                                |  |
|---|--------------------------------|--|
| * | <code>fn_for_warmstart=</code> | Path to the model file (input), saved during training using <code>train</code> via <code>save_fn=</code> . |
|---|--------------------------------|--|

---

Note that the network configuration such as the number of neurons and pooling setting are saved in the model file, and they cannot be changed when training warm-starts. Training parameters such as the regularization parameter and step-sizes need to be specified when warm-starting, and they can be different from before. See `sample/sample-warmstart.sh` to see which parameters should be omitted and which parameter should not at the time of warm-start. Also note that suspending training and resuming training by warm-start would result in a slightly different model from the model trained without interruption. This is mainly because random number generation for mini batch selection is reset at the time of warm-starting.

**Example 8** (`sample/sample-warmstart.sh`). *Train a CNN and save the model after 20 epochs, and warm-start training from the saved model.*

```

conText -l train  save_interval=20 save_fn=output/s-train-p3.mod ...
conText -l train  fn_for_warmstart=output/s-train-p3.mod.ite20 ...

```

*Note that the model filename is a concatenation of `save_fn`, `.ite`, and the number of epochs (20 in this example).*

## 2.5 conText train: AdaDelta

The default training algorithm is mini-batch SGD as described above. Optionally, AdaDelta [7] can also be used.

---

### AdaDelta parameters

---

|   |                          |  |
|---|--------------------------|--|
| * | <code>AdaDelta</code>    | Enable AdaDelta training                         |
| * | <code>[ℓ]rho=</code>     | $\rho$ in Algorithm 1 of [7]. Useful value: 0.95 |
| * | <code>[ℓ]epsilon=</code> | $\epsilon$ in Algorithm 1 of [7].                |

---

When AdaDelta is enabled, the SGD parameters (`step_size` and `momentum`) are ignored. However, the step-size scheduling is in effect, e.g., if `step_size_scheduler=Few`, `step_size_decay=0.1` and `step_size_decay_at=80`, then the amount of weight update by AdaDelta is reduced to one tenth after 80 epochs. It appears that this technique also stabilizes weight updating by AdaDelta as well as SGD and sometimes improves performance. See `sample/sample-adad.sh` for example.

## 2.6 conText train: rmsprop

The default training algorithm is mini-batch SGD as described above. Optionally, *rmsprop* [6] can also be used.

| rmsprop parameters |                                   |   |
|--------------------|-----------------------------------|---|
| *                  | <code>Rmsp</code>                 | Enable rmsprop training                                   |
|                    | <code>[\ell]rmsprop_decay=</code> | Coefficient of exponential moving average. (Default:0.95) |
| *                  | <code>[\ell]step_size=</code>     | Learning rate.  |

When rmsprop is enabled, `momentum` is ignored if specified. The step-size scheduling is in effect if specified.

## 2.7 conText train: layer-by-layer training of single-connection networks

`focus_layers` specifies which layers should be updated during the training. This parameter is valid only on single-connection networks.

|                            |  |
|----------------------------|--|
| <code>focus_layers=</code> | Layers whose weights should be updated. (Default:All)<br>For example, <code>focus_layers=0-1_0-1-2</code> indicates that layer-0 and layer-1 should be updated for the number of epochs specified by <code>num_iterations</code> and designated step-size scheduling first, and then layer-0, 1, and 2 should be updated for <code>num_iterations</code> epochs with the designated step-sized scheduling. This parameter is valid only on single-connection networks. |
|----------------------------|--|

## 2.8 conText predict (applying a CNN to new data)

`conText predict` applies a model saved during training to new data and write prediction values to a file.

**Prediction file (binary output)** By default, `predict` writes prediction values to a file in the following format.

| offset | length                | type        | description   |
|--------|-----------------------|-------------|---|
| 0      | 4                     | int         | <i>s</i> : Length of float. Typically 4.  |
| 4      | 4                     | int         | <i>c</i> : Number of classes  |
| 8      | 4                     | int         | <i>d</i> : Number of data points  |
| 12     | $s \times c \times d$ | float array | Prediction values sorted by data points first and classes next. That is, the first <i>c</i> values are for the first data point, and the next <i>c</i> values are for the second data point, and so on. |

**Prediction file (text output)** Optionally, `predict` writes prediction values in the text format, one data point per line. However, this could be very inefficient in space/time when the number of data points or the number of classes is large.

---

|  |   |
|--|---|
| <code>conText predict (to apply a CNN to test data)</code> |   |
| * <code>model_fn=</code>                                   | Path to the model file (input), saved during training using <code>train</code> via <code>save_fn=</code> parameter. Note that the model filename is a concatenation of <code>save_fn</code> and <code>.ite<i>n</i></code> where <i>n</i> is the number of epochs, e.g., <code>model.ite100</code> if <code>save_fn=model</code> . |
| * <code>prediction_fn=</code>                              | Path to the prediction file (output).   |
| * <code>tstname=</code>                                    | Test data name.   |
| * <code>datatype=</code>                                   | Dataset type. <code>sparse</code>   <code>sparse_multi</code> . This must be the same as training.  |
| <code>data_dir=</code>                                     | Directory where data files are.   |
| * <code>x_ext=</code>                                      | Feature filename extension. This must be the same as training.  |
| <code>test_mini_batch_size=</code>                         | Mini-batch size for parallel processing for testing. (Default: 100)   |
| <code>WriteText</code>                                     | Write prediction values in the text format. (Default: binary format)  |
| <code>extension=</code>                                    | If the model uses multi-connection architecture, specify <code>extension=multi</code> .   |
| <code>data_ext<i>i</i>=</code>                             | If <code>datatype=sparse_multi</code> , specify data name extensions <code>data_ext0</code> , <code>data_ext1</code> , ..., similar to training.  |

---

**Example 9** (`sample/sample-save-predict.sh`). *Single-connection example.*

```
conText gpu# train    save_fn=output/s-train-p3.mod save_interval=10 ...

conText gpu# predict  model_fn=output/s-train-p3.mod.ite20 \
                        prediction_fn=output/s-test-p3.pred.bin \
                        datatype=sparse tstname=s-test-p3 data_dir=data \
                        x_ext=.xsmatvar
```

**Example 10** (`sample/sample-multiconn-save-predict.sh`). *Multi-connection example.*

```
conText gpu# predict  model_fn=output/s-train-p2p3.mod.ite20 \
                        prediction_fn=output/s-test-p2p3.pred.bin \
                        extension=multi datatype=sparse_multi \
                        tstname=s-test- data_ext0=p2 data_ext1=p3 \
                        x_ext=.xsmatvar data_dir=data
```

## 2.9 `conText write_features` (write internal features to a file)

`write_features` writes features produced internally in a CNN to a file. At the moment, this function is limited to writing features of *fixed-sized*. By default, the features produced as input to the top layer are written. Input to a middle layer (which takes input from some other layer) can also be requested via the `layer` parameter; however, if it is not a valid layer (e.g., a layer with variable-sized input), it would result in producing a file of size 0.

**Output file format** The output file is an array of floating-point values (`float` in C++ if the `-D` options of the makefile are unchanged), sorted by data points first and vector dimensions next. For example, if the feature vectors are 3-dimensional and there are 2 data points, the output file contains 24 bytes of six values,

$$f[0, 0], f[1, 0], f[2, 0], f[0, 1], f[1, 1], f[2, 1]$$

in this order, where  $f[i, d]$  is the *i*-th component of the vector for the *d*-th data point.

See `sample/sample-write-features.sh` and `sample/sample-multiconn-write-features.sh` for examples.

**NOTE** CONTEXT v2.00 or higher is required for using this function.

---

| <b>Parameters for</b> <code>conText write_features</code> . |   |
|---|---|
| * <code>model_fn=</code>                                    | Path to the model file (input), saved during training using <code>train</code> via <code>save_fn=</code> parameter. Note that the model filename is a concatenation of <code>save_fn</code> and <code>.ite<i>n</i></code> where <i>n</i> is the number of epochs, e.g., <code>model.ite100</code> if <code>save_fn=model</code> . |
| * <code>feature_fn=</code>                                  | Path to the file to which the features will be written (output).  |
| * <code>tstname=</code>                                     | Data name.  |
| * <code>datatype</code>                                     | Dataset type. <code>sparse sparse_multi</code> . This must be the same as training.   |
| * <code>data_dir=</code>                                    | Directory where data files are.   |
| * <code>x_ext=.xsmatvar</code>                              | Region filename extension.  |
| * <code>layer=</code>                                       | Internal features produced as input to this layer will be written to the designated file. Default: <code>top layer</code> . The layer must be a middle or the top layer with fixed-sized input. Otherwise, it would result in either an error or output of size 0.  |
| <code>extension=</code>                                     | This must be the same as training.  |
| <code>data_ext<i>i</i></code>                               | If <code>datatype=sparse_multi</code> , specify data name extensions <code>data_ext0</code> , <code>data_ext1</code> , ..., similar to training.  |

---

## 2.10 GPU-related settings

This section describes GPU-related parameters, which may be useful for speeding up training/testing in some cases. Basic knowledge of GPU and CUDA programming is assumed.

### 2.10.1 Device memory handler

Allocation of GPU device memory (`cudaMalloc`) is relatively expensive and avoidance of frequent device memory allocation sometimes speeds up processing. For this purpose, CONTEXT is equipped with a device memory handler which obtains a large amount of device memory from the system at the beginning and allocates to the callers when needed. The device memory handler is enabled by the `memory_size` parameter as follows:

Usage: `conText gpu#[:memory_size] action param1 param2 ...`

`memory_size` specifies the size of device memory that the device memory handler should pre-allocate; it is in GB ( $2^{30}$  bytes) and can be a fraction. If `memory_size` is omitted or set to a non-positive value, the device memory handler is disabled. If `memory_size` is too small to satisfy a request by the application with the pre-allocated memory, the memory handler calls `cudaMalloc` so that the application will not fail.

Use of the device memory handler is particularly recommended with LSTM, as it often makes it two times faster. In other cases, the amount of speed up may be small.

Note that when device memory is very tight, use of the device memory handler may cause device memory shortage due to fragmentation, and the device memory handler should be disabled in that case.

**Example 11.** *The following example uses GPU#1 and the device memory handler with 4.5GB device memory.*

```
conText 1:4.5 train...
```

### 2.10.2 To control computation on GPU

The following parameters may be useful.



---

| <b>To control GPU parallel processing in <code>conText</code> <code>train</code> and <code>predict</code></b> |   |
|---|---|
| <code>gpu_max_threads=</code>   | Maximum number of threads per block to be generated on GPU. (Default: the maximum of the GPU)   |
| <code>gpu_max_blocks=</code>  | Maximum number of blocks to be generated on on GPU. (Default: the maximum of the GPU)   |
| <code>NoCusparseIndex</code>  | Do not use cusparse for row indexing (Default: use cusparse). When mini-batch size is small, this sometimes speeds up training/testing significantly. |

---

### 3 Semi-supervised one-hot CNN (using `prepText` and `conText`)

*To get started* Try `sample/sample-ss-unsup-tv.sh` or `exp-in-papers/semisup-imdb-unsup-tv.sh`

The purpose of providing the semi-supervised CNN code is to enable the reproduction of the experiments presented in [2]. Therefore, readers are assumed to be familiar with the concepts/terminology used in [2]. Semi-supervised CNN training described below uses the functions in the earlier sections, and readers are assumed to be familiar with at least Section 1 (data preparation) and Section 2.1 (single-connection CNN training).

Semi-supervised learning with CNNs described in [2] is done in two steps:

1. **Training with unlabeled data:** Obtain a region embedding from unlabeled data. This is done by training with unlabeled data for the artificial task of predicting adjacent regions from each region.
2. **Training with labeled data:** Train a supervised CNN with labeled data for the task of interest using the region embedding obtained above to produce additional input to the CNN.

In both steps, training can be done by `conText train`, which is described for use for supervised learning earlier in this document, and a few additional parameters useful/required for the semi-supervised learning setting will be described below.

**NOTE1** `CONTEXT` v2.00 or higher is required.

**NOTE2** To apply the semi-supervised learning process described here to your data, it is recommended to start with one of the sample scripts and modify it for your data.

- Toy examples: `sample/sample-ss-*.sh`
- Real-world data examples: `exp-in-papers/semisup-*.sh`  
To run these scripts, you need to download `unlab_data.tar.gz` and extract it at `exp-in-papers/` so that the directory `exp-in-papers/unlab_data` will be created.

#### 3.1 Training with unlabeled data (tv-embedding learning)

From unlabeled data, we learn a *tv-embedding* ('tv' stands for 'two-view', see [2] for definition) via training for predicting adjacent regions from each region. This is done by first generating training data from unlabeled data using `prepText gen_regions_unsup` (or `gen_regions_parsup`) and then performing training using `conText train`.

### 3.1.1 Training data generation from unlabeled data

The first step of tv-embedding learning is to generate training data from unlabeled data, for the task of predicting adjacent regions (target regions) using each region (feature region). `prepText gen_regions_unsup` is provided for this purpose. Using this function, feature regions are represented by either bag-of-word, bag-of- $n$ -grams, or position-sensitive vectors (sequential). Target regions are represented by either bag-of-word or bag-of- $n$ -gram vectors.

Optionally, `prepText gen_regions_parsup` can be also used, which represents target regions in a partially-supervised fashion. To use this function, first a CNN needs to be trained with labeled data and the trained CNN should be applied to unlabeled data to obtain the *embedded regions* derived from each region (i.e., output of supervised region embedding) using `conText write_embedded`. The embedded regions are considered to be representations of regions with respect to the task objective (e.g., sentiment analysis). The obtained embedded regions are used as input to `gen_regions_parsup` to represent target regions.

The following table describes parameters for `prepText gen_regions_unsup` and `gen_regions_parsup`.

| <b>Parameters for <code>prepText gen_regions_unsup</code> and <code>gen_regions_parsup</code>.</b> |   |
|--|---|
| * <code>input_fn=</code>   | Path to the input token file or the list of token files. If the filename ends with <code>.lst</code> , the file should be the list of token filenames. The input file(s) should contain one document per line, and each document should be tokens delimited by space.   |
| * <code>x_vocab_fn=</code>   | Path to the vocabulary file generated by <code>gen_vocab</code> , used for X (features). To represent feature regions by bag-of- $n$ -gram vectors, include $n$ -grams in the vocabulary.   |
| * <code>x_type=</code>   | Vector representation for X (features). Bow   Seq. (Default: Bow)   |
| * <code>region_fn_stem=</code>   | Pathname stem of the region vector file, target file, and word-mapping file (output). To make the pathnames, the respective extensions will be attached.  |
| * <code>patch_size=</code>   | Feature region size.  |
| * <code>patch_stride=</code>   | Feature region stride. (Default: 1)   |
| * <code>padding=</code>  | Feature padding size. Feature region size minus one is recommended.   |
| * <code>dist=</code>   | Size of adjacent regions (target regions) used to produce Y (target). The same value as <code>patch_size</code> is recommended.   |
| <code>LowerCase</code>   | Convert upper-case to lower-case characters.  |
| <code>UTF8</code>  | Convert UTF8 en dash, em dash, single/double quotes to ascii characters.  |
| <code>MergeLeftRight</code>  | Do not distinguish the target regions on the left and right.  |
| <b>Parameters for <code>prepText gen_regions_unsup</code> only.</b>                                |   |
| * <code>y_vocab_fn=</code>   | Path to the vocabulary file generated by <code>gen_vocab</code> , used for Y (target). To represent the target regions by bag-of- $n$ -gram vectors, include $n$ -grams in the vocabulary.  |
| <b>Parameters for <code>prepText gen_regions_parsup</code> only.</b>                               |   |
| * <code>embed_fn=</code>   | Pathname to the file produced by <code>conText write_embedded</code> .  |
| * <code>f_patch_size=</code>   | Region size that was used by <code>write_embedded</code> to generate the file. It must be no greater than <code>dist</code> .   |
| * <code>f_patch_stride</code>  | Region stride used by <code>write_embedded</code> .   |
| * <code>f_padding=</code>  | Padding size used by <code>write_embedded</code> .  |
| <code>num_top</code>   | $k$ . To produce target (output) of each instance, only the $k$ largest components of embedded regions will be retained, and the rest will be set to zero. Use this only when the components of embedded regions are guaranteed to be non-negative, e.g., when activation is rectifier, sigmoid, and so on. This speeds up tv-embedding training combined with <i>negative sampling</i> . Useful values: 5, 10. |
| <code>scale_y</code>   | Scale the target values so that the maximum value becomes approximately this value. Useful value: 1.  |

To use `gen_regions_parsup`, `conText write_embedded` should be done in advance to apply a CNN to unlabeled data and obtain embedded regions. Since at the moment the sole purpose of this function is to produce input to `gen_regions_parsup`, the format description of the output binary file is omitted. The table below describes parameters for `write_embedded`.

---

| <b>Parameters for <code>conText write_embedded</code>.</b> |  |
|--|--|
| * <code>tstname=</code>                                    | Data name. Region filename without extension. It must not contain a plus sign.   |
| * <code>datatype=sparse</code>                             | Dataset type. Sparse input.  |
| <code>data_dir=</code>                                     | Directory where data files are.  |
| * <code>x_ext=.xsmatvar</code>                             | Region filename extension.   |
| * <code>model_fn=</code>                                   | Path to the model file (input).  |
| * <code>embed_fn=</code>                                   | Path to the file to which embedded regions will be written (output).   |
| <code>num_top=</code>                                      | <i>k</i> : If specified, only the <i>k</i> largest components of the embedded regions are retained in each vector while the rest are set to zero. Use this only when the embedded regions are guaranteed to be non-negative, e.g., when activation is rectifier, sigmoid, and so on. |

---

### 3.1.2 Training of tv-embedding

Using the training data generated from unlabeled data by `gen_regions_unsup` (or `gen_regions_parsup`), tv-embedding training is done using `conText train` in a manner similar to CNN training for supervised learning. Therefore, all the parameters for `conText train` introduced earlier are available; in addition, the following parameters are required/useful.

---

| <b>Additional parameters for <code>conText train</code>: useful (probably only) for tv-embedding training</b> |   |
|---|---|
| <code>zero_Y_ratio=</code>  | <i>r</i> : Sample <i>negative examples</i> so that only <i>r</i> times more negative examples than <i>positive examples</i> are retained. A <i>negative</i> (or <i>positive</i> ) example is a pair of a region vector and a component of a target vector whose value is zero (or non-zero), respectively. Useful values: 5, 10.  |
| <code>zero_Y_weight=</code>   | <i>t</i> : To compute loss, assign <i>t</i> times larger weights to negative examples than positive examples. Useful values: $1/r$ where <code>zero_Y_ratio=r</code> .  |
| * <code>save_lay0_fn=</code>  | Layer-0 is saved to files specified by this parameter; more precisely, filenames are generated by concatenating <code>iten.layer0</code> where <i>n</i> is the number of epochs. Files are generated at the timing designated by <code>save_interval</code> . When trained for tv-embedding learning, layer-0 embodies the region tv-embedding (a function to convert a region to a vector). The file can be later used to initialize a <i>side layer</i> to produce additional input to a CNN trained with labeled data. |
| <code>NoCusparseIndex</code>  | Do not use the cusparse library for indexing rows of sparse matrices. In the typical tv-embedding training setting (as in the sample scripts), this speeds up training a lot.   |

---

Note that for tv-embedding training `x_ext=.xsmat` and `y_ext=.ysmat` must be specified since those are the types of files `prepText gen_regions_unsup` (or `gen_regions_parsup`) produces. We save a tv-embedding (a function to convert a region to a vector) learned from unlabeled data to a file specified by `save_lay0_fn` and use it in the final training with labeled data as described next.

## 3.2 Training with labeled data and tv-embedding

The final step of semi-supervised learning is to train a CNN with labeled data, using a tv-embedding (obtained from unlabeled data and saved to a file via `save_lay0_fn`) to produce additional input to the CNN. In this code, we call the layer that produces additional input a *side layer*. Conceptually, a side layer is a convolution layer, which receives a document represented by a sequence of one-hot vectors as input and produces a vector for each small region of the document; the layer which the side layer is attached to receives the output of the side layer as input, in addition to the one-hot vector representation of the document. In this code, conversion of one-hot vectors to region vectors (which is done by a convolution layer, conceptually) is done by `prepText gen_regions` in advance; thus, input to a side layer needs to be prepared by `gen_regions` in advance as well.

To specify a training parameter for a side layer, we attach a *side layer identifier* at the beginning, in the form of

`lsidei_`

which refers to the  $i$ -th side layer ( $i = 0, 1, \dots$ ) attached to layer- $\ell$ . For example, if layer-0 has two side layers (i.e., layer-0 receives additional input from two side layers) and layer-1 has one side layer, then there are three valid side layer identifiers: `0side0_`, `0side1_`, and `1side0_`.

---

| Additional parameters for <code>conText train</code> : to specify side layers |   |
|---|---|
| <code>lnum_sides=</code>  | Number of side layers attached to layer- $\ell$ . (Default: 0)  |
| <code>lsidei_fn=</code>   | The file to which <code>conText train</code> saved a tv-embedding via <code>save_lay0_fn</code> .<br>The side layer warm-starts from this file. |
| <code>lsidei_dsno=</code>   | The id of the dataset to be used as input to the side layer.  |
| <code>lsidei_Fixed</code>   | Do not update the weights.  |

---

**NOTE** By default, the side layer weights are updated during CNN training. To prevent updating, `[l]Fixed` should be specified with an appropriate layer id  $\ell$ . In contrast, `reText` fixes side layer weights by default.

## 4 LSTM training and testing: `reText`

**To get started** Try `sample/sample-lstm.sh` or `exp-in-papers/sup-lstm-imdb.sh`.

`reText` is an executable for training and testing LSTM networks using GPU.

Usage: `reText gpu#[memory_size] action param1 param2 ...`

**gpu#** The first argument specifies the GPU ID number. Specify `-1` if the default GPU should be used,

**memory\_size** Specify the size of GPU device memory that the device memory handler should pre-allocate; it is in GB ( $2^{30}$  bytes) and can be a fraction. If omitted or set to a non-positive value, the device memory handler is disabled. For LSTM, a *sufficiently large value* is recommended as it often makes processing *two times faster*. See also 2.10.1.

**action** It must be one of the following:

- `train`: Perform training. Optionally, test the model on development/test data as training proceeds.
- `predict`: Apply the trained model to new data.

*param1, param2, ...* Parameters can be specified via either command line arguments or a parameter file or a combination of both. If an argument starts with an @, the string following @ is regarded as a pathname to a parameter file. See `sample/param/sample.param` for an example of the parameter file, which is for `conText` but the format is the same. Use of a parameter file is sometimes convenient as parameters can be organized and commented.

**Example 12.** Use parameters in `lstm-sample.param` and set `random_seed=1`. Use the default GPU and pre-allocate 4GB device memory.

```
reText -l:4 train random_seed=1 @lstm-sample.param
```

## 4.1 reText train: training and testing LSTM

**To get started** Try `sample/sample-lstm.sh` or `exp-in-papers/sup-lstm-imdb.sh`.

Since `reText train` shares most of the parameters with `conText train`, it is recommended to read Section 2.1 (single-connection CNN) first and come back to this section. Also note that similar to CNN training, data preparation by `prepText gen_vocab` and `gen_regions` (Section 1) is required.

With `reText`, data simply flows from layer-0 to layer-1, layer-1 to layer-2, and so on. A more complex structure such as bi-directional LSTM (combining two LSTM layers) is encapsulated into one layer. This is easier to use but less flexible than the “multi-connection” extension of `conText`. Hidden layers of all types can take input of variable-sized. Input to the top layer must be fixed-sized. There are two types of parameters, *layer parameters* which specifies parameters for each layer, and *network parameters* which are not specific to layers.

### 4.1.1 Network parameters

Network parameters for `reText train` are the same as those for `conText train` described in Section 2.1.1.

### 4.1.2 Layer parameters

Unlike `conText train` (which provides only one type of hidden layer that does everything), `reText` provides distinct types of hidden layer such as an ‘LSTM layer’ and a ‘pooling layer’. Different types of layers require different layer parameters.

Layer parameters can be specified for a specific layer by attaching the layer id  $\ell$  ( $\ell = 0, 1, \dots, \text{top}_-$ ) in front of the keywords, e.g., `0step_size=0.1` for layer-0 or `top_reg_L2=1e-4` for the top layer. Without a layer id, the parameter is regarded as a default parameter and applied to all the applicable layers; i.e., the parameter with a layer id supersedes the one without the id. However, as an exception, omission of the layer id is not allowed for `layer_type` (layer type) and `layer_fn` (layer initialization file).

| Layer parameters.  |   |   |
|--------------------|---|---|
| $\ell$ layer_type= | LstmF   | Forward LSTM (left to right).   |
|                    | LstmB   | Backward LSTM (right to left).  |
|                    | Lstm2   | Bi-directional LSTM (forward LSTM + backward LSTM).   |
|                    | Weight  | Weighting and activation (computation of $\sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ ).<br>If one vector per data point is given as input, a <code>Weight</code> layer serves as a fully-connected layer. |
|                    | Patch   | Concatenate vectors of each region, where regions are specified via size, stride and padding. Dense input only.   |
|                    | Pooling   | Perform pooling and response normalization.   |
| $\ell$ layer_fn=   | Lay   | A layer trained by <code>conText train</code> .   |
|                    | Optionally, a layer can be initialized by a <i>layer file</i> , which was generated via <code>save_layer_fn</code> during training. Used for pre-training and semi-supervised learning. |   |

|                                |   |
|--------------------------------|---|
| <code>[ℓ]save_layer_fn=</code> | <i>x</i> : Layer- <i>ℓ</i> is saved to a <i>layer file</i> , <code>x.iter.layerℓ</code> . where <i>n</i> is the number of epochs at the time of saving. |
| <code>[ℓ]dsno=</code>          | Dataset id. Equivalent to <code>[ℓ]dataset_no</code> in “To use multiple datasets” in Section 2.2. Ignored by Pooling and Patch layers. (Default:0)     |

Optionally, an LSTM layer and a Weight layer can perform pooling and response normalization in this order, after the main job (i.e., doing LSTM or weighting+activation). `reText train` shares most of layer parameters with `conText train`. The table below shows, for each layer type, groups of parameters in Section 2.1.2 that should be referred to.

| Parameter groups in Section 2.1.2        | Lstm{F B} | Weight | Patch | Pooling |
|--|-----------|--------|-------|---------|
| Weighting and activation                 |           | ×      |       |         |
| Weight initialization and regularization | ×         | ×      |       |         |
| SGD/AdaDelta/RMSprop parameters          | ×         | ×      |       |         |
| Pooling, response normalization          | ×         | ×      |       | ×       |
| Region size, stride, and padding         |           |        | ×     |         |

In addition, the following parameters are available for the LSTM layers.

| Layer parameters for an LstmF, LstmB, or Lstm2 layer |   |
|--|---|
| * <code>[ℓ]nodes=</code>                             | Number of LSTM units, or the dimensionality of output vectors.  |
| <code>[ℓ]chop_size=</code>                           | <i>s</i> : during training, each document will be chopped into segments of maximum size <i>s</i> . Recommended values: 50, 100. |
| <code>[ℓ]NoGate_i</code>                             | Do not use the input gate.  |
| <code>[ℓ]NoGate_o</code>                             | Do not use the output gate.   |
| <code>[ℓ]NoGate_f</code>                             | Do not use the forget gate. Not recommended.  |
| <code>[ℓ]bprop_max=</code>                           | Backpropagation goes back only this far. Not recommended. Provided only for baseline experiments. Use chopping instead.         |

**NOTE** An Lstm2 layer is a combination of two layers (LstmF and LstmB). Use layer ID `ℓ-0_` (e.g., “0-0\_”) for LstmF and `ℓ-1_` (e.g., “0-1\_”) for LstmB to specify layer-specific parameters.

### 4.1.3 Output

See Section 2.1.3.

## 4.2 reText predict (applying a saved model to new data)

`reText predict` applies a model saved during training to new data and writes prediction values to a file. The interface is the same as `conText predict` in Section 2.8, except for the `extention` parameter being irrelevant. See `sample/sample-lstm-save-predict.sh` for an example.

## 5 Semi-supervised LSTM (using `prepText` and `reText`)

*To get started* Try `exp-in-papers/semi-lstm2-imdb.sh`.

The purpose of providing the semi-supervised LSTM code is to enable the reproduction of the experiments presented in [3]. Therefore, readers are assumed to be familiar with the concepts/terminology used in [3]. Semi-supervised LSTM training described below uses the functions in the earlier sections, and readers are assumed to be familiar with at least Section 1 (data preparation) and Section 4 (LSTM training and testing).

Similar to semi-supervised one-hot CNN, semi-supervised LSTM training described in [3] is done in two steps:

1. **Training with unlabeled data:** Train an LSTM with unlabeled data for the artificial task of predicting an adjacent future region at every time step.
2. **Training with labeled data:** Train an LSTM with labeled data for the task of interest using the LSTM trained above to produce additional input.

In both steps, training can be done by `reText train`, described in Section 4, and a few additional parameters useful/required for semi-supervised LSTM training will be described below.

**NOTE1** `CONTEXT` v3.00 or higher is required.

**NOTE2** To apply the semi-supervised learning process described here to your data, it is recommended to start with one of the sample scripts `exp-in-papers/semi-*.sh` and modify it for your data.

- To run these scripts, you need to download `unlab_data.tar.gz` and extract it at `exp-in-papers/` so that the directory `exp-in-papers/unlab_data` will be created.
- To use the provided tv-embedding files, you need to download `lstm_output.tar.gz` and extract it at `exp-in-papers/` so that the directory `exp-in-papers/provided_output` will be created.

### 5.1 LSTM training with unlabeled data (LSTM tv-embedding learning)

From unlabeled data, we learn a *tv-embedding* ('tv' stands for 'two-view', see [2] for definition) via training for predicting an adjacent future region at each time step. This is done by first generating training data from unlabeled data using `prepText gen_regions_unsup` and then performing training using `reText train`.

#### 5.1.1 Training data generation from unlabeled data

The first step of tv-embedding learning is to generate training data from unlabeled data, for the task of predicting an adjacent future region (target region) at each time step. `prepText gen_regions_unsup`, introduced in Section 3.1.1, is used for this purpose with additional parameters below.

---

| <b>Additional parameters for <code>prepText gen_regions_unsup</code>.</b> |  |
|---|--|
| <code>RightOnly</code>  | Use this for training a forward LSTM (left to right) so that the region to the right (future) of the current time step is regarded as a target region. |
| <code>LeftOnly</code>   | Use this for training a backward LSTM (right to left) so that the region to the left (future) of the current time step is regarded as a target region. |

---

### 5.1.2 Training of LSTM tv-embedding

Using the training data generated from unlabeled data by `gen_regions_unsup`, LSTM tv-embedding training is done using `reText train` in a manner similar to supervised LSTM training. Therefore, all the parameters for `reText train` introduced earlier are available; in addition, the following parameters are required/useful.

---

| Additional parameters for <code>reText train</code> : useful (probably only) for tv-embedding training |                    |
|--|--------------------|
| <code>zero_Y_ratio=</code>   | See Section 3.1.2. |
| <code>zero_Y_weight=</code>  | See Section 3.1.2. |
| <code>NoCusparseIndex</code>   | See Section 3.1.2. |

---

We save an LSTM layer, which embodies a tv-embedding (a function to convert a region to a vector) learned from unlabeled data, to a file specified by `save_layer_fn` and use it in the final training with labeled data as described next.

## 5.2 Training with labeled data and LSTM tv-embedding

The final step of semi-supervised LSTM training is to train LSTM with labeled data, using tv-embeddings (obtained from unlabeled data and saved to files via `save_layer_fn`) to produce additional input. As in Section 3, we call the layer that produces additional input a *side layer*.

To specify layer parameters for a side layer, we attach a *side layer identifier* at the beginning, in the form of

`0side $i$ _`

which refers to the  $i$ -th side layer ( $i = 0, 1, \dots$ ) attached to layer-0, e.g., `0side0_` and `0side1_`. Note that unlike semi-supervised CNN where side layers can be attached to any layer, semi-supervised LSTM allows only layer-0 to have side layers; therefore, the side layer identifier always starts with `0side`.

---

| Additional parameters for <code>reText train</code> : to specify side layers |  |
|--|--|
| <code>num_sides=</code>  | Number of side layers attached to layer-0.                           |
| <code>ℓside<math>i</math>_UpdateSide</code>                                  | Update the weights of side layers. (Default: fix side layer weights) |

---

**Example 13.** *Training with two LSTM tv-embeddings (forward and backward). The side layer weights are fixed.*

```
reText -1:5 train layers=2 num_sides=2 \  
0layer_type=Lstm2 0chop_size=50 \  
0side0_layer_type=LstmF 0side0_layer_fn=lstmF.itel0.layer0 \  
0side1_layer_type=LstmB 0side1_layer_fn=lstmB.itel0.layer0 \  
1layer_type=Pooling 1num_pooling=1 1pooling_type=Max \  
:
```

**Example 14.** *Training with two LSTM tv-embeddings (forward and backward) and one CNN tv-embedding. The weights of all three side layers are fixed. The third side layer (CNN tv-embedding) takes input from dataset-1, and other bottom layers take input from dataset-0 (default). See also “To use multiple datasets” in Section 2.2.*



```

reText -l:5 train layers=2 num_sides=3 \
0layer_type=Lstm2 0chop_size=50 \
0side0_layer_type=LstmF 0side0_layer_fn=lstmF.itel0.layer0 \
0side1_layer_type=LstmB 0side1_layer_fn=lstmB.itel0.layer0 \
0side2_layer_type=Lay 0side2_layer_fn=cnn.itel0.layer0 \
0side2_Fixed 0side2_dsno=1 \
1layer_type=Pooling 1num_pooling=1 1pooling_type=Max \
:

```

## 6 Examples of other usage

There are other network architectures that can be explored using `conText` and `reText`. Sample scripts are provided for some of them.

### 6.1 One-hot CNN with a convolution layer above a convolution layer

`sample/sample-cc.sh` trains a CNN with a convolution layer above a convolution layer.

### 6.2 One-hot CNN with a fully-connected layer above a convolution layer

`sample/sample-cf.sh` trains a CNN with a fully-connected layer above a convolution layer.

### 6.3 Fully-connected neural networks

`sample/sample-1nn.sh` and `sample/sample-2nn.sh` train a fully-connected neural network of one layer and two layers, respectively, using bag-of-word vectors as input.

### 6.4 Linear models

`sample/sample-linear.sh` trains a linear model using bag-of-word vectors as input.

### 6.5 Word-vector CNN

The experimental results in [1, 2] and <http://riejohnson.com/software/cnn-perf.html> lead to the following recommendation on CNN configurations:

- If your training data is huge, then try supervised one-hot CNN (Section 2).
- Else if you have good unlabeled data (large and in-domain), then try semi-supervised one-hot CNN (Section 3).
- Else if your training set is a small number of short texts (snippets/sentences), try *word-vector CNN* with public general-purpose word vectors (this section). Unlabeled data, even if it's general, helps a lot, especially in this setting.
- Else, you may need to try more than one thing above.

A word embedding layer (or word table lookup, the first layer of word-vector CNN) can be implemented by a convolution layer of region size 1 where the columns of the weights serve as word vectors. To use pre-trained word vectors (such as publicly available word vectors) as input to CNN (and LSTM also), they need to be transformed to a *weight file* so that `conText train` can use them for initializing the weights.

---

| <b>Parameters for <code>conText adapt_word_vectors</code> to convert word vectors to a weight file.</b> |   |
|---|---|
| * <code>wordvec_bin_fn</code>   | Path to a binary file containing word vectors and words in the <code>word2vec</code> binary format (input). Either this or <code>wordvec_txt_fn</code> below is required. Not both. Note that <code>conText</code> reads this file in the endian-sensitive fashion. To share word vectors among the systems with different endianness, the text format below can be used instead. |
| * <code>wordvec_txt_fn</code>   | Path to a text file containing word vectors and words in the <code>word2vec</code> text format (input). Either this or <code>wordvec_bin_fn</code> above is required. Not both.   |
| * <code>word_map_fn=</code>   | Path to the word-mapping file ( <code>*.xtext</code> ) generated by <code>prepText gen_regions</code> (input). Word vectors will be sorted in the order of this file so that the dimensions of the resulting weights will correctly correspond to the dimensions of the region vectors generated by <code>gen_regions</code> .  |
| * <code>weight_fn=</code><br><code>rand_param=</code>   | Path to the weight file to be generated (output).<br><i>x</i> : scale of initialization. If the word-mapping file ( <code>word_map_fn</code> ) contains words for which word vectors are <i>not</i> given, the word vectors for these unknown words will be randomly set by Gaussian distribution with zero mean with standard deviation <i>x</i> . (Default:0)                   |
| <code>random_seed=</code><br><code>IgnoreDupWords</code>  | Seed of random number generation.<br>Ignore it if there are duplicated words associated with word vectors. If this is not turned on, the process will be terminated on the detection of duplicated words.   |

---



---

| <b>Parameters for <code>conText train</code> (and the <code>Weight</code> layer of <code>reText train</code>) to initialize weights by a weight file.</b> |  |
|---|--|
| <code>[ℓ]weight_fn</code>   | Path to a weight file that contains initial weights. |

---

See `sample/sample-ss-wordvec.sh` for an example.

## 7 Hints and Tips

### 7.1 Miscellaneous tips

Here are things to consider for obtaining good performance.

- **seq-CNN vs. bow-CNN.** With `Bow` turned on when generating region files, region files for bow-CNN is generated, which represents regions by bag-of-word vectors; otherwise, seq-CNN is trained which represents regions by concatenation of one-hot vectors. Roughly speaking, if on your data, linear models with bag-of-*n*-gram vectors outperform linear models with bag-of-word vectors, then seq-CNN with a small region size (e.g., 3, 4, or 5) is likely to be effective on that data. Otherwise, bow-CNN with a large region size (e.g., 10, 15, or 20) is likely to be effective on that data. In [1], an example application of the former is sentiment classification and the latter is news topic classification.
- **Step-size scheduling.** In the experiments in [1], the number of epochs was fixed to 100 and the step-size was always reduced by multiplying 0.1 after 80 epoch (`num_iterations=100 step_size_scheduling=Few step_size_decay=0.1 step_size_decay_at=80`). This technique (reducing the step-size once or twice towards the end of training), also used in [4], often improves performance while it never hurts it.

- **Initial step-size.** It is important to set the initial step-size `step_size` to a good value. To search for a good value on the development data, typically, we tested 0.5, 0.25, 0.1, 0.05,  $\dots$ , while fixing momentum to 0.9 and mini-batch size to 100 (`momentum=0.9 mini_batch_size=100`) and step-size scheduling to the setting above.
- **Number of epochs.** The number of training data points used in [1] was around 25K. 50–100 epochs seemed appropriate for this size of data. If your training data is much larger, fewer epochs may suffice since there is more data to go through.
- **Case shifting.** Be consistent on case shifting. If you generated a vocabulary file with the `LowerCase` switch turned on, then generate regions files with the `LowerCase` switch on.
- The scripts in the `exp-in-papers` directory are for reproducing the experiments reported in [1, 2, 3], and they provide realistic examples on the real-world data (movie reviews, product reviews, and news). Note that the sample scripts at the `sample` directory are only for showing how to use `CONTEXT`, and their parameter settings are not intended for the optimum performance.

## 7.2 GPU memory consumption during training (`conText train`, `reText train`)

In a typical environment, GPU memory (*device memory*) is relatively small compared with CPU memory (*host memory*). When `conText` (or `reText`) runs out of GPU memory, it terminates with an error message saying “out of memory”, e.g.,

```
!cuda error!: (Detected in AzPmem::alloc extra)
AzCuda::throwIfError: AzPmat::reform_noinit 8.2e+08
cudaGetErrorString returned out of memory
```

The following describes, *with v3 or higher*, what are loaded into GPU memory during training by `conText train` (or `reText train`) and therefore what parameters affect the consumption of GPU memory.

- **Weights.** A larger number of weight vectors (specified by `nodes`) consumes more GPU memory. A larger number of target classes also consumes more GPU memory since the number of weight vectors in the top layer is equal to the number of classes.
- **Data in a mini batch** and work areas for processing a mini batch. A larger mini batch size (`mini_batch_size`) consumes more GPU memory. Also, a larger `num_pooling` increases the size of feature vectors used in the layer above, which are in GPU memory.

With v3 or higher, *consumption of GPU device memory does not depend on training data size or test data size.*

## 7.3 CPU memory consumption during training (`conText train`, `reText train`)

**Training data** Training data is read from the region file and target file and stays in the CPU memory during training. If training data is too large to fit in your CPU memory, divide the data into several *batches*. By doing so, only one of the batches is loaded into CPU memory; thus, for example, having  $k$  batches reduces CPU memory consumption to  $1/k$ . Note, however, that each batch must be a good sample of the feature/class distributions; therefore, division into batches should be done randomly. See Section 2.3 and `sample/sample-batch.sh` for how to use batches.

**Test data** If optional testing is done, test data also stays in the CPU memory. `conText train` (or `reText train`) is not meant for very large test data. `conText predict` (or `reText predict`) should be used for very large test data, if necessary, multiple times. That is, if your test data is large, avoid testing while training (which can be done by specifying `NoTest` and omitting `tstname`), save a trained model to a file (using `save_fn` and `save_interval`), and apply the saved model to the test data using `conText predict` (or `reText predict`).

## 7.4 When to make training data batches (`conText train`, `reText train`)

The batch option should be used if the amount of training data is so large that

- it does not fit in the CPU memory, or
- one epoch (going through the entire training data) takes very long and you want to save/test the models or reduce the step-size before reaching the end of the training data.

Note that each batch must be a good sample of the feature/class distributions; therefore, *division into batches must be done randomly*. See Section 2.3 and `sample/sample-batch.sh` for how to use batches.

## References

- [1] Rie Johnson and Tong Zhang. Effective use of word order for text categorization with convolutional neural networks. In *Proceedings of NAACL HLT*, 2015.
- [2] Rie Johnson and Tong Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In *Proceedings of NIPS*, 2015.
- [3] Rie Johnson and Tong Zhang. Supervised and semi-supervised text categorization using LSTM for region embeddings. In *Proceedings of ICML*, 2016.
- [4] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *Proceedings of NIPS*, 2012.
- [5] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(Jun):1929–1958, 2014.
- [6] Tijman Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4, 2012.
- [7] Matthew D. Zeiler. ADADELTA: An adaptive learning rate method. *arXiv:1212.5701*, 2012.